

1 **Discordance between different bioinformatic methods for identifying**  
2 **resistance genes from short-read genomic data, with a focus on**  
3 ***Escherichia coli***  
4

5 **1.1 Author names**

6 *Timothy J Davies*<sup>a, b</sup>, *Jeremy Swan*<sup>a, b</sup>, *Anna E Sheppard*<sup>a, b</sup>, *Hayleah Pickford*<sup>a, b</sup>,  
7 *Samuel Lipworth*<sup>a, b</sup>, *Manal AbuOun*<sup>c</sup>, *Matthew Ellington*<sup>b, e</sup>, *Philip W Fowler*<sup>a</sup>,  
8 *Susan Hopkins*<sup>b, e</sup>, *Katie L Hopkins*<sup>b, f</sup>, *Derrick W Crook*<sup>a, b, d</sup>, *Tim EA Peto*<sup>a, b, d</sup>,  
9 *Muna F Anjum*<sup>c</sup>, *A Sarah Walker*<sup>a, b (\*)</sup>, *Nicole Stoesser*<sup>a, b, d (\*)</sup>.

10 \* contribution considered equal

11 **1.2 Affiliation**

- 12 a) *Nuffield Department of Medicine, Oxford University, Oxford, United Kingdom*  
13 b) *National Institute for Health Research (NIHR) Health Protection Research Unit*  
14 *on Healthcare Associated Infections and Antimicrobial Resistance at*  
15 *University of Oxford, UK*  
16 c) *Bacteriology, Animal and Plant Health Agency, Surrey UK*  
17 d) *Oxford University Hospitals NHS Foundation Trust, Oxford, UK*  
18 e) *Antimicrobial Resistance and Healthcare Associated Infections (AMRHAI)*  
19 *Division, UK Health Security Agency, London, UK*  
20 f) *HCAI, Fungal, AMR, AMU and Sepsis Division, UK Health Security Agency,*  
21 *London, UK*  
22

23 **1.3 Corresponding author**

24 *Dr Timothy Davies, [timothy.davies@ndm.ox.ac.uk](mailto:timothy.davies@ndm.ox.ac.uk)*

25 *Alternate corresponding author:*

26 *Dr Nicole Stoesser, [nicole.stoesser@ndm.ox.ac.uk](mailto:nicole.stoesser@ndm.ox.ac.uk)*

27

28 **1.4 Keyword**

29 *Antimicrobial resistance genotyping, genomics, Escherichia coli, resistance*  
30 *prediction*  
31

32 **1.5 Repositories:**

33 Sequencing data are available at the following NCBI BioProject accession number:  
34 PRJNA540750.

35 **2. Abstract**

36 Several bioinformatics genotyping algorithms are now commonly used to  
37 characterise antimicrobial resistance (AMR) gene profiles in whole genome  
38 sequencing (WGS) data, with a view to understanding AMR epidemiology and  
39 developing resistance prediction workflows using WGS in clinical settings. Accurately  
40 evaluating AMR in Enterobacterales, particularly *Escherichia coli*, is of major  
41 importance, because this is a common pathogen. However, robust comparisons of  
42 different genotyping approaches on relevant simulated and large real-life WGS  
43 datasets are lacking. Here, we used both simulated datasets and a large set of real

44 *E. coli* WGS data (n=1818 isolates) to systematically investigate genotyping methods  
45 in greater detail.

46

47 Simulated constructs and real sequences were processed using four different  
48 bioinformatic programs (ABRicate, ARIBA, KmerResistance, and SRST2, run with  
49 the ResFinder database) and their outputs compared. For simulations tests where  
50 3,092 AMR gene variants were inserted into random sequence constructs,  
51 KmerResistance was correct for all 3,092 simulations, ABRicate for 3,082 (99.7%),  
52 ARIBA for 2,927 (94.7%) and SRST2 for 2,120 (68.6%). For simulations tests where  
53 two closely related gene variants were inserted into random sequence constructs,  
54 ABRicate identified the correct alleles in 11,382/46,279 (25%) of simulations, ARIBA  
55 in 2494/46,279 (5%), SRST in 2539/46,279 (5%) and KmerResistance in  
56 38,826/46,279 (84%). In real data, across all methods, 1392/1818 (76%) isolates  
57 had discrepant allele calls for at least one gene.

58

59 Our evaluations revealed poor performance in scenarios that would be expected to  
60 be challenging (e.g. identification of AMR genes at <10x coverage, discriminating  
61 between closely related AMR gene sequences), but also identified systematic  
62 sequence classification (i.e. naming) errors even in straightforward circumstances,  
63 which contributed to 1081/3092 (35%) errors in our most simple simulations and at  
64 least 2530/4321 (59%) discrepancies in real data. Further, many of the remaining  
65 discrepancies were likely “artefactual” with reporting cut-off differences accounting  
66 for at least 1430/4321 (33%) discrepant. Comparing outputs generated by running  
67 multiple algorithms on the same dataset can help identify and resolve these  
68 artefacts, but ideally new and more robust genotyping algorithms are needed.

69

### 70 **3. Impact statement**

71 Whole-genome sequencing is widely used for studying the epidemiology of  
72 antimicrobial resistance (AMR) genes in bacteria; however, there is some concern  
73 that outputs are highly dependent on the bioinformatics methods used. This work  
74 evaluates these concerns in detail by comparing four different, commonly used AMR  
75 gene typing methods using large simulated and real datasets. The results highlight  
76 performance issues for most methods in at least one of several simulated and real-  
77 life scenarios. However most discrepancies between methods were due to  
78 differential labelling of the same sequences related to the assumptions made  
79 regarding the underlying structure of the reference resistance gene database (i.e.  
80 that resistance genes can be easily classified in well-defined groups). This study  
81 represents a major advance in quantifying and evaluating the nature of  
82 discrepancies between outputs of different AMR typing algorithms, with relevance for  
83 historic and future work using these algorithms. Some of the discrepancies can be  
84 resolved by choosing methods with fewer assumptions about the reference AMR  
85 gene database and manually resolving outputs generated using multiple programs.  
86 However, ideally new and better methods are needed.

87

88

## 89 4. Introduction

90 Whole genome sequencing (WGS) has become a major tool for characterising the  
91 epidemiology of bacterial antimicrobial resistance (AMR) genes, representing a  
92 potentially highly discriminatory, non-targeted approach with significant advantages  
93 over other more targeted molecular techniques(1). In addition, WGS-based antibiotic  
94 susceptibility prediction has been successfully implemented as part of diagnostic and  
95 treatment workflows for *Mycobacterium tuberculosis*(2). Accurate WGS-based  
96 profiling of complete AMR gene content and prediction of susceptibility phenotypes  
97 would represent an attractive option for other commonly encountered clinical  
98 bacterial pathogens, such as Enterobacterales, including *Escherichia coli*.  
99

100 Several key components are required for WGS-based AMR genotyping and  
101 predictions of susceptibility phenotype, including a robust AMR gene reference  
102 catalogue linking each genetic mechanism/sequence with a given phenotype, and  
103 accurate AMR gene identification and classification algorithms. Several catalogues  
104 and bioinformatics algorithms are now available(3-9), but only limited comparative  
105 evaluation of their outputs has been undertaken. The genetic mechanisms  
106 underpinning AMR in Enterobacterales and some other bacteria (e.g. *Pseudomonas*  
107 *aeruginosa*) are much more complex than those in *M. tuberculosis*, and whilst some  
108 studies suggest that WGS-based genotyping holds promise for AMR gene  
109 characterisation and the prediction of antimicrobial susceptibility for several different  
110 Enterobacterales species(10-12), the limited reproducibility and reliability of such  
111 methods in a blinded, head-to-head analysis across nine bioinformatics teams has  
112 been recently highlighted(13). However, this study was small (n=10 sequencing  
113 datasets, n=7 isolates), encountered a limited set of typing discrepancies, and used  
114 highly selected samples, meaning the impact of these issues on larger, real-world  
115 datasets remains unclear.  
116

117 We therefore used simulations and three large, independent and diverse *E. coli*  
118 sequencing datasets to investigate the robustness and reproducibility of four widely-  
119 used WGS-based AMR genotyping methods (ABRicate, ARIBA, KmerResistance,  
120 and SRST2) at scale, investigating any encountered discrepancies.  
121

## 122 5. Methods

### 123 *AMR gene identification methods*

124 We evaluated the impact of different bioinformatics tools using the same AMR gene  
125 catalogue, namely the ResFinder database (v.29/10/2019). At the time the study was  
126 designed (March 2018), to be included bioinformatics tools had to: (i) have publicly  
127 available code, (ii) run on local computing architecture without major modification,  
128 (iii) accept different AMR gene databases to ensure broad and long-term typing  
129 usability, and (iv) have a command line interface that could enable batch processing  
130 of large numbers of samples (**Table S1**).  
131

132 We identified four publicly available bioinformatic tools that met these criteria and  
133 used distinct AMR gene identification approaches: ABRicate(14) (which searches for  
134 AMR genes in assemblies using BLASTn), SRST2(7) (which maps reads directly  
135 onto the formatted AMR gene database using Bowtie 2), ARIBA(6) (which combines

136 these two approaches, first mapping reads to the AMR gene database using  
137 minimap, and then creating local assemblies of the mapped reads using Fermi-lite)  
138 and KmerResistance(8) (which analyses shared k-mers between the query  
139 sequences and reference sequences in the AMR gene database) (**Fig.S1**). To mimic  
140 broad usability, each program was run using default parameters. For ABRicate,  
141 assemblies were first produced using SPAdes(15) run with default parameters.

142

143 *Simulated data: single and multiple allele identification, and low coverage scenarios*

144 Prior to evaluating real data, we considered the accuracy of each method in  
145 identifying known AMR gene alleles “inserted” into simulated flanking sequence  
146 constructs. For this, each AMR gene variant in the ResFinder database (n=3,092)  
147 was flanked by 1kb of random sequence (using Numpy v1.16.4(16) and combined  
148 using BioPython(17) v1.74) and reads simulated at 40x coverage using ART (details  
149 and rationale in Supplementary Methods, **Fig.1, S2**). Other ART parameters were:  
150 error profile=“HISEQ2500”, mean DNA fragment length (standard deviation)=480bp  
151 (150bp), and read length=151bp. Each bioinformatic method was then tested to see  
152 if it could correctly identify the AMR gene variant, using default parameters.

153

154 We also considered two *a priori* scenarios that are thought to affect AMR  
155 genotyping(18), namely a *multiple allele* scenario in which multiple closely  
156 genetically related alleles (see below) of a given AMR gene were present, and a *low*  
157 *quality* scenario reflected by low sequencing coverage. For the *multiple allele*  
158 scenario we excluded target AMR gene variants that were incorrectly identified  
159 individually by any method (see Results), and then calculated pairwise nucleotide  
160 similarity between all remaining AMR gene variants. To do this, each remaining AMR  
161 gene variant was split into 31-mers, which were then compared with 31-mer sets  
162 from every other non-excluded AMR gene variant using pairwise Jaccard’s similarity  
163 indices. AMR gene variant pairs were defined as similar if they shared any 31-mer,  
164 resulting in a total of 46,279 possible similar AMR gene variant pairs (**Fig.S3-S5**).

165

166 For the *low coverage* scenario, reads were simulated from 176 *bla*<sub>TEM</sub> gene-  
167 containing constructs at coverage depths ranging from 1x to 50x using ART  
168 (n=176\*50=8,800 simulations), reflecting total *bla*<sub>TEM</sub> diversity present in the  
169 ResFinder database at the time of simulation. Each construct contained a random  
170 perfect reference *bla*<sub>TEM</sub> variant flanked by 1kb of random sequence on each side  
171 produced using Numpy/BioPython as above. Simulated reads were then processed  
172 by each genotyping method using default settings and the identified variants were  
173 compared with the known *bla*<sub>TEM</sub> variants present in each construct. The measure of  
174 performance for this scenario was the proportion of *bla*<sub>TEM</sub> variants correctly  
175 identified by each method at each coverage level.

176

177 *Real data: Isolate selection*

178 To evaluate performance on real data, we then studied a total of 1,818 *E. coli*  
179 isolates comprising three different WGS datasets in order to reflect different strain-  
180 level and AMR gene diversity: (i) 984 sequentially collected bloodstream infection  
181 isolates at Oxford University Hospitals (OUH) NHS Foundation Trust(19) (“Oxford  
182 dataset”); (ii) 497 animal commensal *E. coli* isolates donated by the UK Animal and  
183 Plant Health Agency (APHA)(20) (“APHA dataset”), and (iii) 337 *E. coli* isolates  
184 collected by UK Health Security Agency’s (UKHSA) Antimicrobial Resistance and

185 Healthcare Associated Infections (AMRHAI) Reference Unit, which investigates  
186 isolates enriched for rare or important resistance genotypes encountered in the UK  
187 (sequenced for this study, “UKHSA dataset”).  
188

189 Isolates were re-cultured from frozen stocks stored in nutrient broth plus 10%  
190 glycerol at -80°C. DNA was extracted using the QuickGene DNA Tissue Kit S  
191 (Kurabo Industries, Japan) as per manufacturer’s instructions, with an additional  
192 mechanical lysis step (FastPrep, MP Biomedicals, USA) immediately following  
193 chemical lysis. A combination of standard Illumina and in-house protocols were used  
194 to produce multiplexed paired-end libraries, which were sequenced on an Illumina  
195 HiSeq 2500, generating 151bp paired-end reads. High quality sequences were de-  
196 novo assembled using Velvet(21) as previously described(22). *In silico* Achtman(23)  
197 multi-locus sequence types (MLST) types were defined using ARIBA(6).  
198

199 While this work does not attempt to predict resistance from WGS data, each isolate  
200 had linked AST (summarized in **Table S2, Fig.S6**), which we have included as the  
201 complexity of resistance genotype identification is associated with the phenotype.  
202 Isolates had complete AST data available for: ampicillin, ceftazidime and one other  
203 3rd generation cephalosporin (cefotaxime for the animal commensal isolates,  
204 ceftriaxone for all others), gentamicin, ciprofloxacin, and co-trimoxazole.  
205

206 We compared AMR genotypes reported for each isolate by each method, stratified  
207 by antibiotic class to which resistance was conferred as specified in the ResFinder  
208 database, namely: beta-lactams, aminoglycosides, quinolones, trimethoprim, and  
209 sulphonamides. Discrepancies were classified according to which of the four  
210 bioinformatics methods agreed (**Fig.S7**). The cause of discrepancy was investigated  
211 for all beta-lactam resistance genotypes, because these antibiotics are most  
212 commonly used for clinical *E. coli* infections, and then for discrepancy patterns  
213 occurring in >1.5% (n=27) of isolates for the other classes.

## 214 **6. Results**

### 215 ***Simulated scenarios***

#### 216 *Accurate identification of single AMR gene variants in simulated sequence* 217 *constructs*

218 For the 3,092 AMR gene variants in the ResFinder database, all four genotyping  
219 methods correctly identified those inserted into random sequence contexts in 2,011  
220 (63.5%) cases. KmerResistance was correct for all 3,092 simulations, ABRicate for  
221 3,082 (99.7%), ARIBA for 2,927 (94.7%) and SRST2 for 2,120 (68.6%) (**Fig.2**). For  
222 SRST2, most errors were due to its approach of pre-clustering reference sequences  
223 into sub-families by sequence identity prior to genotyping, thereby essentially  
224 excluding *a priori* the possibility of identifying alleles that were not selected as the  
225 representative for these sub-family clusters. This error is explained in more detail  
226 below as it also affected genotyping in real isolate sequences.  
227

#### 228 *Impact of the presence of multiple closely related alleles on genotyping calls*

229 The multiple allele simulation caused significant problems for assembly-based  
230 algorithms, with ABRicate reporting fragmented/incomplete alleles for 32,194/46,279  
231 (70%) simulations and ARIBA reporting no alleles meeting its assembly quality  
232 requirements for 32,987/46,279 (71%) simulations. SRST2, as expected, found only

233 a single allele in most (33077/46,279 (71%)) cases (**Table 1**), as dictated by its  
234 clustering parameters. ABRicate managed to identify both alleles correctly in the  
235 absence of incorrect calls in 11,382/46,279 (25%) of simulations, whereas ARIBA  
236 and SRST2 only managed to correctly reconstruct both members of the pair in the  
237 absence of correct calls in 2,494/46,279 (5%) and 2,539/46,279 (5%) cases  
238 respectively (Table 1). Of the four programs, KmerResistance performed the best,  
239 identifying both alleles correctly without additional erroneous calls in 38,826/46,279  
240 (84%) of cases. Unsurprisingly all four programs were most likely to make  
241 erroneous genotyping calls as the simulated pairs of alleles became more closely  
242 related (**Fig.S8**).

243

#### 244 *Impact of sequencing depth on genotyping calls*

245 KmerResistance was able to identify *bla*<sub>TEM</sub> alleles at lower coverage than any of the  
246 other methods (**Fig.1**). Above 15x depth of coverage for the gene, all methods  
247 correctly identified *bla*<sub>TEM</sub> alleles in simulated constructs in > 95% of cases (**Fig.1**).  
248 All methods were able to identify all of the *bla*<sub>TEM</sub> alleles correctly at least once, but  
249 examples existed for all methods where the allele was correctly identified at low  
250 coverage, but then mis-classified at higher coverage. In general, ABRicate and  
251 SRST2, while requiring greater sequencing depth to correctly identify *bla*<sub>TEM</sub> alleles  
252 initially were more accurate at higher coverage depths, making erroneous calls for  
253 only 1/176 (0.6%) and 0/176 (0%) of *bla*<sub>TEM</sub> alleles at depths >20x. In contrast, for  
254 >20x coverage ARIBA and KmerResistance made erroneous allele calls for 23/176  
255 (13%) and 6/176 (3%) *bla*<sub>TEM</sub> variants respectively. Above 40x coverage ABRicate  
256 was incorrect for one (0.6%), ARIBA for four (2%), KmerResistance for one (0.6%),  
257 and SRST2 for zero (0%) simulated *bla*<sub>TEM</sub> alleles.

258

#### 259 **Real data**

##### 260 *E. coli* isolate diversity, antimicrobial susceptibility phenotypes and antimicrobial 261 resistance genotypes

262 The 1,818 isolates were diverse, representing >260 multi-locus sequence types  
263 (STs), which were differentially distributed among the datasets. For example,  
264 although ST131 was the most common (207/1818 (11%) isolates), this was largely  
265 due to the fact it was by far the most common in the UKHSA dataset (74/337 (22%)  
266 isolates). In the Oxford dataset, it was only the second most common ST (123/984  
267 (13%) isolates) after ST73 (161/984 (16%)) isolates) and it was rare in the APHA  
268 isolates (10/497 isolates (2%)).

269

270 Correspondingly, the set also contained a broad range of resistance genes, but the  
271 exact number was dependant on the method of search. For legibility, we have  
272 included results as reported by ABRicate as this is the most conceptually simple and  
273 interrogatable approach.. The commonest AMR-associated sequence identified was  
274 *mdfA*. This is known to be universal in *E. coli*, and correspondingly was identified in  
275 all 1,818 isolates in the dataset. There were no other ubiquitous AMR genes;  
276 however, several were common across datasets, with *bla*<sub>TEM</sub>, *aadA*, *sul*, *tet*, and *dfr*  
277 genes occurring in >40% of the isolates. As expected, more UKHSA isolates  
278 contained extended-spectrum beta-lactamase (54/337 vs 94/1481) and  
279 carbapenemase (18/337 vs 1/1481) genes ( $p < 0.001$ ). Aside from *bla*<sub>TEM</sub>, other  
280 beta-lactamases were rare among the APHA dataset. Outside of beta-lactam-

281 associated AMR genes, the Oxford dataset had the lowest proportion of other AMR  
282 genes for all the different gene families encountered in this study.

283

#### 284 *Genotyping discrepancies*

285 10,487 different genes (N=15,588 different alleles) were identified in the 1818  
286 isolates by the four methods. 1,392/1,818 (76%) isolates had discrepancies across  
287 the four bioinformatics methods for at least one gene. At the gene-level, aside from  
288 for *tet*, *aadA* and *cat* genes, the performance of the bioinformatic tools was similar  
289 (**Fig.3, panel a**), with tools reporting each gene in the approximately same  
290 proportion of isolates (within +/-2%). With regards to the three outliers, ABRicate  
291 reported *tet* and *aadA* genes in 19% and 10% more isolates respectively than the  
292 other three tools, and ABRicate and KmerResistance reported *cat* genes in 5% more  
293 isolates than ARIBA and SRST2. By contrast, the alleles reported by each tool were  
294 often discrepant, with alleles of some genes (e.g. *blaSHV*, *blaCMY*) consistently  
295 being differentially reported (**Fig.3, panel b**). Consequently, pairwise agreement  
296 between any two different tools was less than 59% (N=1,065 isolates, **Fig.3, panel**  
297 **c**). While unsupported genotype reports (i.e. where the output of one tool was not  
298 supported by any other) were common for all tools (**Fig.4**), KmerResistance reported  
299 fewer unsupported genotypes than the other three tools ( $p < 0.001$ ).

300

#### 301 *Causes of genotyping discrepancy*

302 At least 2,530/4,321 (59%) of allele-level discrepancies were due to programs  
303 naming the same underlying sequence differently (annotation differences). We  
304 identified three major causes of differences through investigation of discrepantly  
305 reported genes: (i) difficulty distinguishing between optimal matches among alleles  
306 with nested sequences (N=1,737 genes); (ii) spurious identification of additional  
307 alleles due to reads being multiply mapped to distant variants of the same allelic  
308 family (N=547 genes); and (iii) tools choosing different optimal matches based on  
309 DNA sequence alignment when the database only contains one sequence per  
310 protein (N=197) (**Fig.5**). These issues occurred alone in 1,944/2,530 (77%)  
311 discrepantly reported genes, and or in combination in 586/2,530 (23%) cases. In  
312 isolation these errors typically caused only a single method to be discordant, but  
313 when combined resulted in more complex patterns of discrepancy and could make  
314 all four methods disagree with one another. In addition to annotation, ABRicate's  
315 more relaxed requirement for complete gene coverage (which aims to mitigate  
316 assembly errors) caused at least 1,430/4,321 (33%) allele-level discrepancies.  
317 Discrepancies less easily classified as (but likely related to) annotation/cut-offs did  
318 occur, but only affected 381/10487 (4%) of reported genotypes.

319

#### 320 *Annotation-related discrepancies*

321 The most common type of annotation error (N=1,737 genes) was the result of tools  
322 struggling to choose optimal matches where the database contained nested  
323 sequences. One such example of this (N=24) was caused by the sequences for two  
324 different *dfrA7* alleles in the October 2019 Resfinder database, *dfrA7\_1\_AB161450*  
325 and *dfrA7\_5\_AJ419170*. The shorter of the two (*dfrA7\_1\_AB161450*, 474 base pairs  
326 long) aligns almost perfectly (percentage identity = 99%, 1 single nucleotide gap)  
327 with the first 473 bases of *dfrA7\_5\_AJ419170*. ARIBA, KmerResistance and SRST2,  
328 which look for the best identity sequence matches, all report the sample contains a  
329 perfect match for *dfrA7\_1\_AB161450*. By contrast ABRicate, which uses BLAST to

330 identify optimal sequences, reports the sample contains a near perfect match to  
331 *dfrA7\_5\_AJ419170*, as with this being a longer match it is more statistically  
332 significant. Similar errors occurred for several other genes, including *sul*, *tet*, *aph(6)*,  
333 and *aac(3)*.  
334

335 The second most common annotation discrepancy (N=547 genes) represented tools  
336 reporting multiple alleles due to reads mapping to two or more distant variants of the  
337 same allelic family. An example observed was ARIBA and SRST2 reporting multiple  
338 *bla<sub>SHV</sub>* alleles. In this instance, ARIBA and SRST2 identified a primary perfect allele  
339 and a second allele with a lower quality match. These multiple matches however  
340 were likely spurious, with <10 reads mapping individually to each allele, no clear  
341 heterozygosity observed in read pileups, and no fragmentation in assembly graphs.  
342 This is the result of a byproduct of how mapping methods identify optimal matches.  
343 Both ARIBA and SRST2 map reads to each sequence in the database, and then  
344 compare “closely related” sequences to decide which mapping is optimal. Defining  
345 “closely related” however is not straightforward (**Fig.S9**). Reads mapping to more  
346 than one set of “closely related” sequences can result in tools finding multiple gene  
347 variants when the isolate only had one gene original  
348

349 The final common annotation discrepancy (N=197 genes) was due to allele reporting  
350 based on which sequence in the database had the optimal DNA alignment with the  
351 target resistance gene. Although resistance gene nomenclature is largely based on  
352 protein sequence, but resistance gene databases mostly only catalogue one  
353 nucleotide sequence linked to an associated protein sequence. Variant alleles with  
354 synonymous mutations fail to perfectly match any element, and may have an  
355 alternate optimal DNA match. We observed this on 9 occasions where ABRicate,  
356 KmerResistance and SRST2 identified imperfect nucleotide-level matches to  
357 *aph(3'')-lb\_2\_AF024602* and ARIBA identified an imperfect match to *aph(3'')-*  
358 *lb\_4\_AF313472*. However, the sequence they were matching to in the SPAdes and  
359 ARIBA assembly was a 100% identity and coverage protein match to *aph(3'')-*  
360 *lb\_5\_AF321551*.  
361

#### 362 *Non-annotation related discrepancies*

363 In addition to annotation discrepancies that were caused by bioinformatics  
364 algorithms, genotyping calls were also affected by partial/low coverage of AMR gene  
365 targets and assembly fragmentation, consistent with the results from simulations. For  
366 some of these, such as the 1,430 cut-off related discrepancies occurring for *tet*, *mfs*,  
367 *aadA*, and *cat* genes, each program identified the same section of sequence, making  
368 it clear that the different programs had different thresholds for reporting, other  
369 situations were less clear. To investigate this in detail, we examined beta-lactamase  
370 matches which were either partial/low coverage or occurred across fragmented  
371 assemblies.  
372

373 Partial/low coverage beta-lactamase genes were discrepantly found in 39 isolates  
374 (**Fig.S10**), particularly affecting *bla<sub>TEM</sub>*-like gene calls (29/39 cases). KmerResistance  
375 reported the presence of a beta-lactamase gene in all 39 of these discrepant cases,  
376 with calls supported to a varying degree by the other algorithms. However, in all but  
377 four cases, KmerResistance reported that the depth of the gene was less than 5x.  
378 For the four cases where the gene was present at greater than 5x depth as called by



379 KmerResistance, three (present at depth >100x) were omitted from ARIBA reports  
380 as ARIBA assemblies contained mis-sense mutations and the final one (present at  
381 depth 17x) also failed to assemble for ABRicate.

382

383 Assembly fragmentation affected ABRicate and ARIBA beta-lactam resistance gene  
384 calls in 24 cases, with 16 of these likely to be due to the presence of multiple closely  
385 related beta-lactamase alleles affecting assembly integrity. The possibility of  
386 heterozygous alleles was indicated by the ARIBA flag  
387 “variants\_suggest\_collapsed\_repeat”, and the SRST2 “minor allele frequency value”  
388 was high (>20%). KmerResistance reported two related alleles in 12/16 cases, one  
389 with high depth, percentage identity and coverage, and one much less accurately.  
390 This likely reflects KmerResistance’s winner-takes-all strategy, where matching  
391 unique k-mers to reference alleles are counted, and the reference allele with the  
392 most matches is then also assigned all reads with non-unique kmer-matches. This  
393 then leaves only reads with unique k-mers matching any closely related secondary  
394 allele, resulting in poor depth and coverage metrics.

395

## 396 7. Discussion

397 We evaluated the impact of bioinformatics approaches to AMR genotyping in *E. coli*  
398 for four commonly used methods and a widely used AMR gene database  
399 (ResFinder). Using >50,000 simulations and comparing >1,800 sequences sampled  
400 across human and animal reservoirs, thereby capturing common and rare AMR  
401 genotypes, we highlight that whilst currently available, widely-used genotyping  
402 methods are useful, their outputs should be carefully considered in light of our  
403 findings. Commonly postulated causes of discrepancy, such as low quality  
404 sequencing data, appeared to play little role. Instead, discrepancies were primarily  
405 artefactual, occurring because of different approaches in representing the complexity  
406 of the reference AMR gene database. Inconsistent labelling of gene variants will also  
407 affect the reliability of any catalogue-based methods for phenotypic prediction from  
408 WGS-based AMR genotypes. Specifically, predicting phenotype based on the  
409 presence of specific allelic variants will be problematic without a reliable method of  
410 identification.

411

412 Our work agrees with previous findings by Doyle *et al.* on a small and selected  
413 dataset(13); however, we utilised large simulated and real-life datasets to identify  
414 these significant genotyping discrepancies between methods, and also characterized  
415 the underlying reasons for these discrepancies. We found most discrepancies were  
416 largely due to annotation differences, i.e. each method identified the same  
417 consensus sequence but then named them differently. Further, many of these  
418 discrepancies are caused by implicit and frequently incorrect assumptions about  
419 database structure and AMR gene diversity, namely: that AMR genes can be  
420 classified in well-defined families using genetic identity, that different approaches to  
421 deciding best-matching alleles are equivalent, and that isolates will usually not  
422 harbour highly genetically related variants of the same AMR gene. However,  
423 nomenclature and family structure amongst AMR genes relevant to Enterobacterales  
424 is complicated, with highly diverse genotypes (and sometimes phenotypes) being  
425 assigned similar family names (e.g. *bla*<sub>CTX-M</sub>, *bla*<sub>OXA</sub>) and single SNPs in some cases  
426 leading to different resistance phenotypes (e.g. *bla*<sub>TEM-1</sub> (Genbank: AY458016.1) -

427 beta-lactamase inhibitor susceptible i.e. susceptible to amoxicillin-clavulanate, *bla*<sub>TEM-</sub>  
428 <sub>30</sub> (Genbank: AJ437107.1) - beta-lactamase inhibitor resistant i.e. resistant to  
429 amoxicillin-clavulanate). Given this, it is not surprising that we found methods that  
430 make fewer assumptions (e.g. KmerResistance) to be more robust. Based on our  
431 findings accurate resistance genotyping may require the use of multiple different  
432 methods to cross-check results, and a clear understanding of the specific  
433 assumptions underlying the methods used, before conclusions about allele presence  
434 are drawn. The alternative is the development of new algorithms that cope better  
435 with underlying AMR gene diversity in these organisms.

436  
437 One of the key strengths of this analysis was its combined use of both simulations  
438 and real world data. By using simulations, we were able to benchmark methods  
439 against a known truth, which is impossible to do with real-world data. Previous  
440 studies using only real-world data have attempted to overcome the absence of  
441 complete knowledge of the underlying genotype by using phenotypic data as a  
442 reference standard; however genotype-phenotype correlations remain poorly  
443 defined(10, 19). By subsequently using a large sequencing dataset of isolates  
444 obtained across niches, we were then able to assess the extent of discrepancies in  
445 real-life, replicating the problems observed in simulated data.

446  
447 A limitation of this work is that we chose not to evaluate the impact of database  
448 choice, and this will represent future work. Currently, as has been highlighted  
449 previously(24), there are discrepancies between the AMR databases in common  
450 use, with each having a slightly different scope and in some cases differential names  
451 for different AMR gene variants (e.g. *strA* vs *aph(6)-Ia* or *aphD*, and *strB* versus  
452 *aph(6)-Id*). Comparing databases would have therefore added significant further  
453 complexity whilst limiting the generalisability of findings. A further limitation stemming  
454 from our fixed choice of database is that we have not analysed any methods where  
455 the bioinformatic method and database are intertwined (e.g. ResFinder/PointFinder  
456 or RGI). As the interaction between tool and database was the cause of many  
457 issues, it is possible that methods that are database-specific will perform better.  
458 However, the drawbacks of these combined resources are their inflexibility, again  
459 limiting generalisability. A further limitation was that these genotyping algorithms  
460 were compared using an older version of the ResFinder database – the most up to  
461 date when this work was originally planned. Since this time, 70 sequences have  
462 been added, 2 sequences modified and 2 sequences deleted (See supplementary  
463 data). We opted not to re-perform the analysis due to its manual nature and that as  
464 most of the discrepancies relate to underlying principles behind the algorithms rather  
465 than the specific implementation. Finally, we have focused our evaluation on *E. coli*,  
466 but it is likely that these issues will also more widely affect AMR genotyping,  
467 particularly of similar species with complex genotypes.

468  
469 While WGS-based approaches are attractive for both characterizing AMR gene  
470 epidemiology and representing a subsequent tool for resistance prediction, this work  
471 highlights the need for caution when interpreting resistance genotypes reported by  
472 even widely used bioinformatics methods. Before WGS-based approaches can be  
473 considered reliable for use in *E. coli* (and likely other Enterobacterales), particularly  
474 for clinical decision making or replacing phenotypic data to determine

475 epidemiological trends, database standardisation, the development of novel  
476 genotyping approaches, and improved validation and evaluation will be required.  
477

## 478 **8. Author statements**

### 479 **8.1 Authors and contributors**

480 TJD, NS, AES, ASW, DWC and TEAP conceptualised the study. TD, NS, ASW, AES  
481 and MFA decided the methodology. NS, ASW, MFA, AES, DWC and TEAP  
482 supervised the project. NS, MA, MFA, MJE, KH and SH acquired and curated the  
483 data used in this study. TJD and JSW constructed software pipelines to analyse  
484 sequencing data using each of the bioinformatic tools. TJD and ASW investigated  
485 the data. TJD performed the formal analysis. NS, AES, SL, HP, AES and TEAP  
486 assisted with interpreting the cause and impact of discrepancies. TJD and NS wrote  
487 the original draft. TJD, NS, AES, PWF, TEAP and ASW assisted with data  
488 visualisation. All authors were involved in the review and editing process.

### 489 **8.2 Conflicts of interest**

490 The authors have no conflicts of interest to declare.

### 491 **8.3 Funding information**

492 The study was funded by the National Institute for Health Research Health  
493 Protection Research Unit (NIHR HPRU) in Healthcare Associated Infections and  
494 Antimicrobial Resistance at Oxford University in partnership with Public Health  
495 England (PHE) [NIHR200915]. DWC, TEAP, PWF and ASW are supported by the  
496 NIHR Oxford Biomedical Research Centre. The report presents independent  
497 research funded by the National Institute for Health Research. The views expressed  
498 in this publication are those of the authors and not necessarily those of the NHS, the  
499 National Institute for Health Research, the Department of Health or Public Health  
500 England. NS is an Oxford Martin Fellow and an NIHR Oxford BRC Senior Fellow.  
501 ASW is an NIHR Senior Investigator.  
502

### 503 **8.4 Ethical approval**

504 Not applicable.  
505

### 506 **8.5 Acknowledgements**

507 We are grateful to the microbiology laboratory teams at the John Radcliffe Hospital,  
508 Oxford, the Animal and Plant Health Agency, and UK Health Security Agency.

## 509 **9. References**

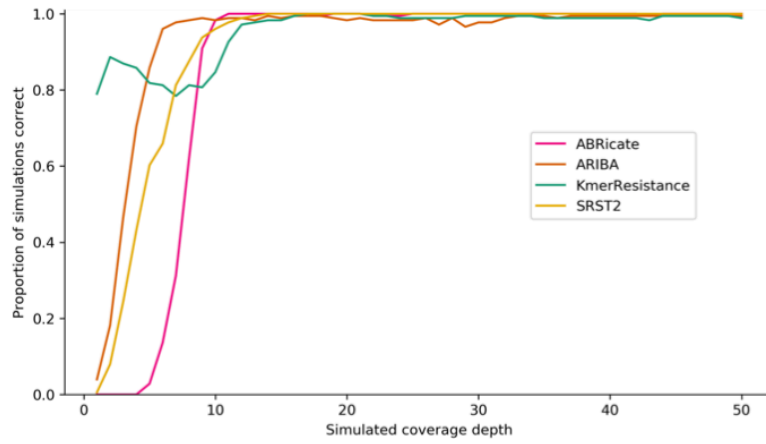
- 510 1. Quainoo S, Coolen JPM, van Hijum S, Huynen MA, Melchers WJG, van  
511 Schaik W, et al. Whole-Genome Sequencing of Bacterial Pathogens: the Future of  
512 Nosocomial Outbreak Analysis. *Clin Microbiol Rev.* 2017;30(4):1015-63.
- 513 2. Quan TP, Bawa Z, Foster D, Walker T, Del Ojo Elias C, Rathod P, et al.  
514 Evaluation of Whole-Genome Sequencing for Mycobacterial Species Identification  
515 and Drug Susceptibility Testing in a Clinical Setting: a Large-Scale Prospective  
516 Assessment of Performance against Line Probe Assays and Phenotyping. *J Clin  
517 Microbiol.* 2018;56(2).

- 518 3. Bortolaia V, Kaas RS, Ruppe E, Roberts MC, Schwarz S, Cattoir V, et al.  
519 ResFinder 4.0 for predictions of phenotypes from genotypes. *J Antimicrob*  
520 *Chemother.* 2020;75(12):3491-500.
- 521 4. Alcock BP, Raphenya AR, Lau TTY, Tsang KK, Bouchard M, Edalatmand A,  
522 et al. CARD 2020: antibiotic resistance surveillance with the comprehensive antibiotic  
523 resistance database. *Nucleic Acids Res.* 2020;48(D1):D517-d25.
- 524 5. Feldgarden M, Brover V, Haft DH, Prasad AB, Slotta DJ, Tolstoy I, et al.  
525 Validating the AMRFinder Tool and Resistance Gene Database by Using  
526 Antimicrobial Resistance Genotype-Phenotype Correlations in a Collection of  
527 Isolates. *Antimicrob Agents Chemother.* 2019;63(11).
- 528 6. Hunt M, Mather AE, Sánchez-Busó L, Page AJ, Parkhill J, Keane JA, et al.  
529 ARIBA: rapid antimicrobial resistance genotyping directly from sequencing reads.  
530 *Microb Genom.* 2017;3(10):e000131.
- 531 7. Inouye M, Dashnow H, Raven LA, Schultz MB, Pope BJ, Tomita T, et al.  
532 SRST2: Rapid genomic surveillance for public health and hospital microbiology labs.  
533 *Genome medicine.* 2014;6(11):90.
- 534 8. Clausen PT, Zankari E, Aarestrup FM, Lund O. Benchmarking of methods for  
535 identification of antimicrobial resistance genes in bacterial whole genome data. *J*  
536 *Antimicrob Chemother.* 2016;71(9):2484-8.
- 537 9. Zankari E, Allesøe R, Joensen KG, Cavaco LM, Lund O, Aarestrup FM.  
538 PointFinder: a novel web tool for WGS-based detection of antimicrobial resistance  
539 associated with chromosomal point mutations in bacterial pathogens. *J Antimicrob*  
540 *Chemother.* 2017;72(10):2764-8.
- 541 10. Stoesser N, Batty EM, Eyre DW, Morgan M, Wyllie DH, Del Ojo Elias C, et al.  
542 Predicting antimicrobial susceptibilities for *Escherichia coli* and *Klebsiella*  
543 *pneumoniae* isolates using whole genomic sequence data. *J Antimicrob Chemother.*  
544 2013;68(10):2234-44.
- 545 11. Shelburne SA, Kim J, Munita JM, Sahasrabhojane P, Shields RK, Press EG,  
546 et al. Whole-Genome Sequencing Accurately Identifies Resistance to Extended-  
547 Spectrum  $\beta$ -Lactams for Major Gram-Negative Bacterial Pathogens. *Clin Infect Dis.*  
548 2017;65(5):738-45.
- 549 12. Stubberfield E, AbuOun M, Sayers E, O'Connor HM, Card RM, Anjum MF.  
550 Use of whole genome sequencing of commensal *Escherichia coli* in pigs for  
551 antimicrobial resistance surveillance, United Kingdom, 2018. *Euro Surveill.*  
552 2019;24(50).
- 553 13. Doyle RM, O'Sullivan DM, Aller SD, Bruchmann S, Clark T, Coello Pelegrin A,  
554 et al. Discordant bioinformatic predictions of antimicrobial resistance from whole-  
555 genome sequencing data of bacterial isolates: an inter-laboratory study. *Microb*  
556 *Genom.* 2020;6(2).
- 557 14. Seemann T. ABRicate. 2020.
- 558 15. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al.  
559 SPAdes: a new genome assembly algorithm and its applications to single-cell  
560 sequencing. *J Comput Biol.* 2012;19(5):455-77.
- 561 16. Harris CR, Millman KJ, van der Walt SJ, Gommers R, Virtanen P,  
562 Cournapeau D, et al. Array programming with NumPy. *Nature.* 2020;585(7825):357-  
563 62.
- 564 17. Cock PJ, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al.  
565 Biopython: freely available Python tools for computational molecular biology and  
566 bioinformatics. *Bioinformatics.* 2009;25(11):1422-3.

- 567 18. Ellington MJ, Ekelund O, Aarestrup FM, Canton R, Doumith M, Giske C, et al.  
568 The role of whole genome sequencing in antimicrobial susceptibility testing of  
569 bacteria: report from the EUCAST Subcommittee. *Clin Microbiol Infect.* 2017;23(1):2-  
570 22.
- 571 19. Davies TJ, Stoesser N, Sheppard AE, Abuoun M, Fowler P, Swann J, et al.  
572 Reconciling the Potentially Irreconcilable? Genotypic and Phenotypic Amoxicillin-  
573 Clavulanate Resistance in *Escherichia coli*. *Antimicrob Agents Chemother.*  
574 2020;64(6).
- 575 20. AbuOun M, O'Connor HM, Stubberfield EJ, Nunez-Garcia J, Sayers E, Crook  
576 DW, et al. Characterizing Antimicrobial Resistant *Escherichia coli* and Associated  
577 Risk Factors in a Cross-Sectional Study of Pig Farms in Great Britain. *Front*  
578 *Microbiol.* 2020;11:861.
- 579 21. Zerbino DR. Using the Velvet de novo assembler for short-read sequencing  
580 technologies. *Curr Protoc Bioinformatics.* 2010;Chapter 11:Unit 11 5.
- 581 22. Stoesser N, Sheppard AE, Peirano G, Anson LW, Pankhurst L, Sebra R, et al.  
582 Genomic epidemiology of global *Klebsiella pneumoniae* carbapenemase (KPC)-  
583 producing *Escherichia coli*. *Sci Rep.* 2017;7(1):5917.
- 584 23. Wirth T, Falush D, Lan R, Colles F, Mensa P, Wieler LH, et al. Sex and  
585 virulence in *Escherichia coli*: an evolutionary perspective. *Mol Microbiol.*  
586 2006;60(5):1136-51.
- 587 24. McArthur AG, Tsang KK. Antimicrobial resistance surveillance in the genomic  
588 age. *Ann N Y Acad Sci.* 2017;1388(1):78-91.

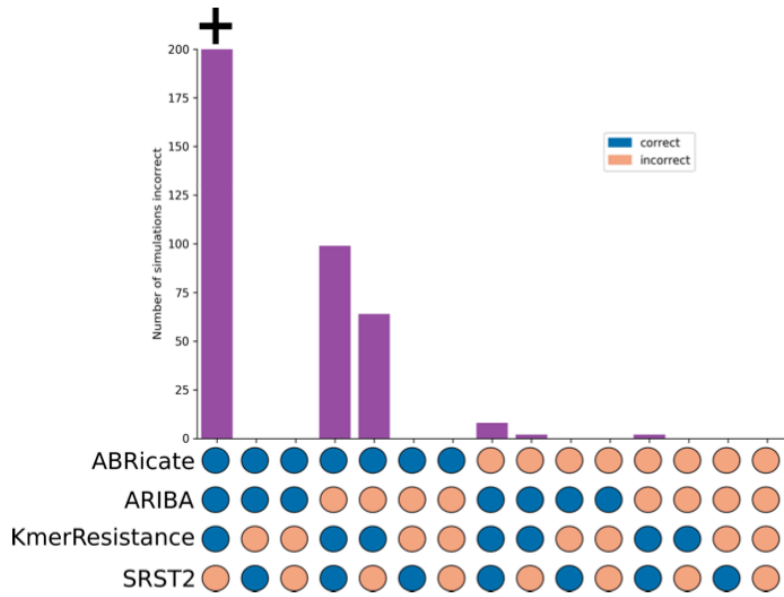
589 **10. Figures and tables**

590 **Figure 1. Proportion of correct genotype calls for single AMR gene variants in**  
591 **simulated constructs by coverage depth and bioinformatics method.**



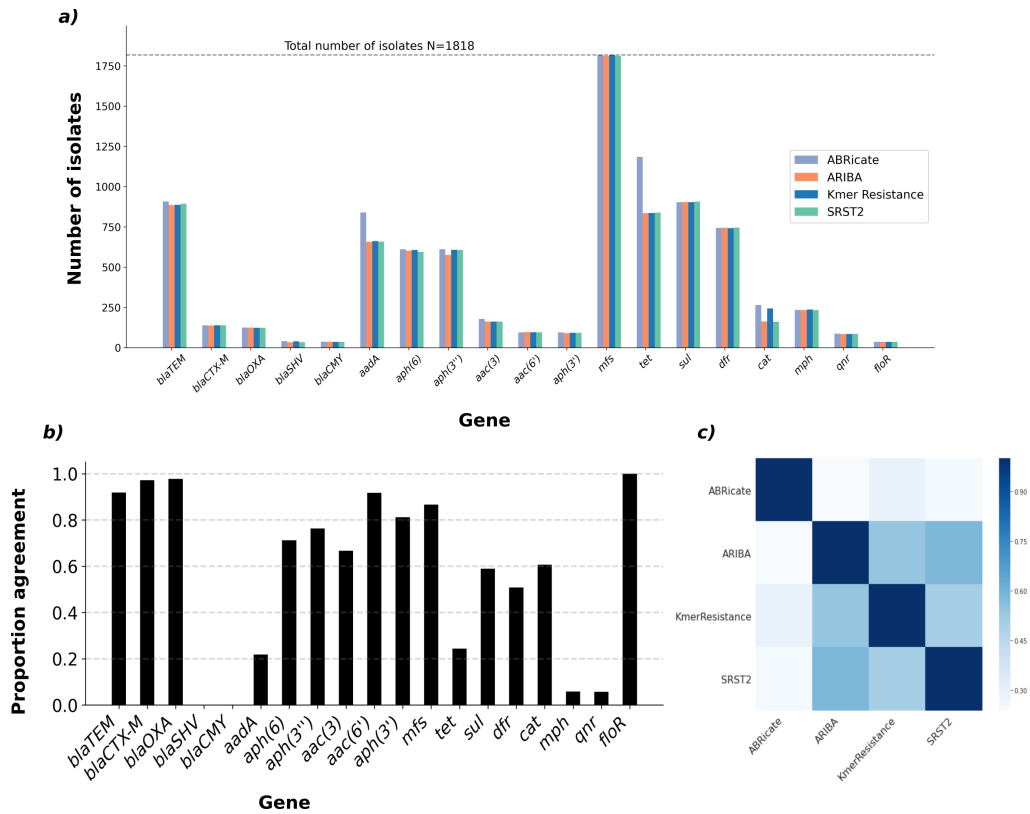
592

593 **Figure 2. Identification of known single AMR gene variants in simulated**  
594 **contexts by bioinformatic method.** Note only cases where one or more methods  
595 were incorrect are shown (n=1,081). “+” denotes the case where total SRST2-only  
596 errors=906, but are truncated to 200 to make other errors visible. blue = method  
597 correct for these simulations, orange = method incorrect.  
598



599  
600  
601

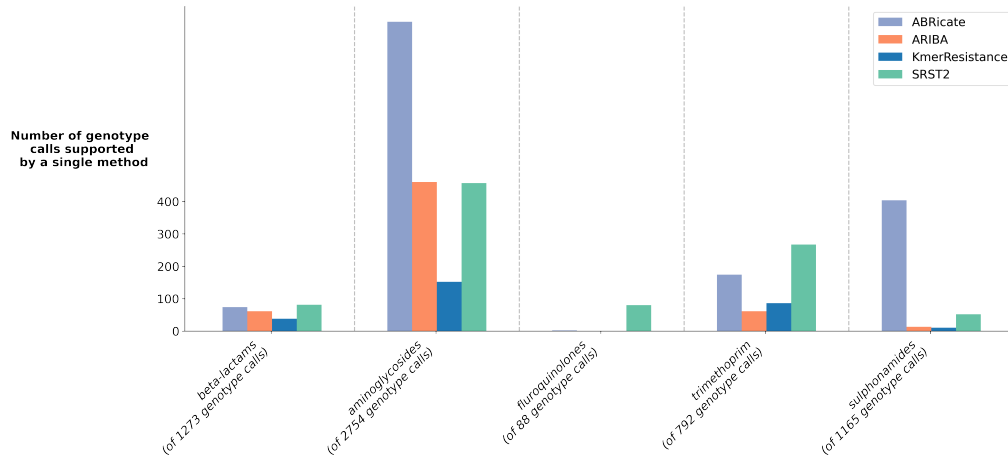
602 **Figure 3. Gene identification concordance vs allele identification concordance.**  
 603 a) The number of isolates containing at least one allele of the name gene families (x-  
 604 axis) stratified by method. b) The proportion of times a given gene was identified  
 605 concordantly by all four methods. c) Pairwise agreement between the different  
 606 methods across all isolates.  
 607



608  
 609  
 610



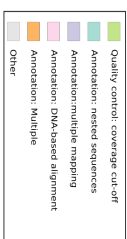
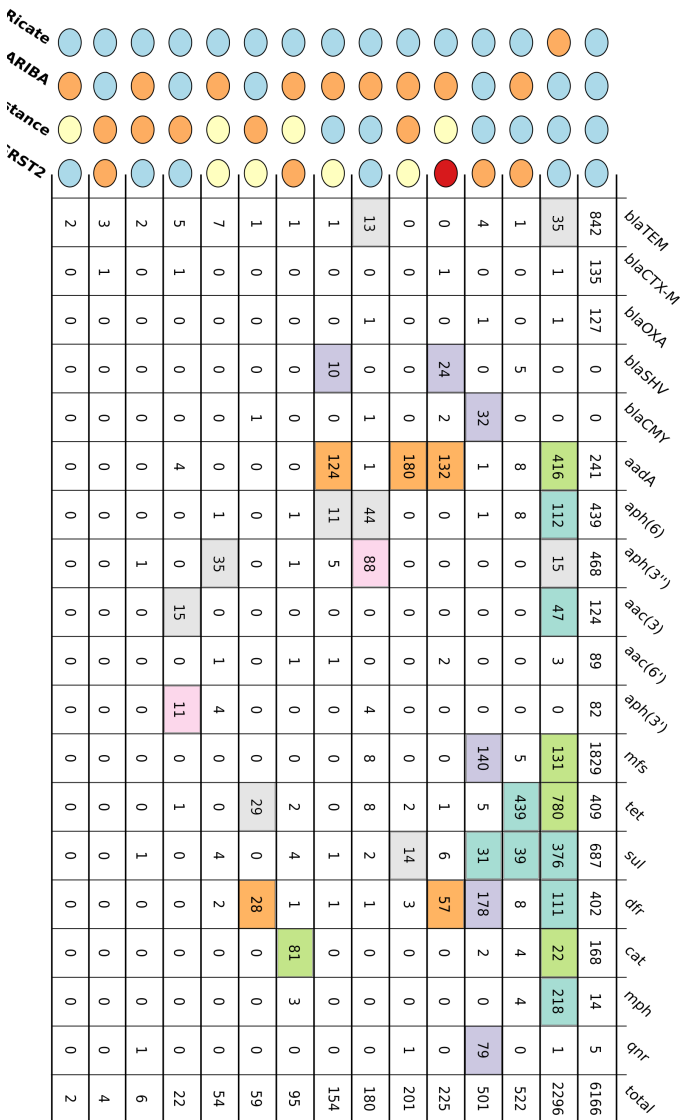
611 **Figure 4. Genotype calls produced by a single method only, stratified by**  
612 **antibiotic class.**  
613



614  
615  
616

617  
618  
619  
620  
621  
622

**Figure 5. Genotyping agreement across all four bioinformatics algorithms, stratified by gene.**  
Colours on the left indicate which methods agreed with one another, with circles with the same colour indicating agreement.  
Colours in the main panel of the figure were used to identify the cause of the discrepancy, as denoted in the figure key. Cells (in the figure) were coloured if > 90% of isolates were caused by a given discrepancy. Cells with <10 isolates were not investigated.



624 **Table 1. Performance of genotyping methods in evaluating simulated**  
625 **constructs with two related allelic variants.** Percentage reported out of a total of  
626 46,279 simulations performed for each method.  
627

Genotyping call	Number of calls (%)			
	ABRicate	ARIBA	KmerResistance	SRST2
No correct calls	17,145 (37%)	36,150 (78%)	489 (1%)	9,898 (21%)
One correct call but additional incorrect calls	2,419 (5%)	2 (0%)	1,452 (3%)	152 (0%)
One correct call, no incorrect calls	15,333 (33%)	7,634 (17%)	2,203 (5%)	33,077 (71%)
Two correct calls, but additional incorrect calls	0 (0%)	1 (0%)	3,309 (7%)	613 (1%)
Two correct calls, no incorrect calls	11,382 (25%)	2494 (5%)	33826 (84%)	2539 (5%)

628